

(19)

(11) **EP 1 599 867 B1**

(12)

EUROPEAN PATENT SPECIFICATION

(45) Date of publication and mention
of the grant of the patent:
13.02.2008 Bulletin 2008/07

(51) Int Cl.:
G10L 15/26 (2006.01)

(21) Application number: **04716133.6**

(86) International application number:
PCT/US2004/006112

(22) Date of filing: **01.03.2004**

(87) International publication number:
WO 2004/079720 (16.09.2004 Gazette 2004/38)

(54) IMPROVING THE TRANSCRIPTION ACCURACY OF SPEECH RECOGNITION SOFTWARE

**VERBESSERUNG DER TRANSKRIPTIONSGENAUIGKEIT VON
SPRACHERKENNUNGSSOFTWARE**

AMELIORATION DE LA PRECISION DE TRANSCRIPTION DE LA RECONNAISSANCE VOCALE

(84) Designated Contracting States:
**AT BE BG CH CY CZ DE DK EE ES FI FR GB GR
HU IE IT LI LU MC NL PL PT RO SE SI SK TR**

(30) Priority: **01.03.2003 US 451024**

(43) Date of publication of application:
30.11.2005 Bulletin 2005/48

(73) Proprietors:
• **Colfman, Robert E.**
Millville,
New Jersey 08332 (US)
• **Bender, Frederick J.**
Rosenhayen, NJ 08352 (US)

(72) Inventors:

- **Colfman, Robert E.**
Millville,
New Jersey 08332 (US)
- **Bender, Frederick J.**
Rosenhayen, NJ 08352 (US)

(74) Representative: **Beresford, Keith Denis Lewis et al**
Beresford & Co.,
16 High Holborn
London WC1V 6BX (GB)

(56) References cited:
WO-A-01/26093 **WO-A-01/69905**
GB-A- 2 345 783 **US-A- 5 758 322**

EP 1 599 867 B1

Note: Within nine months from the publication of the mention of the grant of the European patent, any person may give notice to the European Patent Office of opposition to the European patent granted. Notice of opposition shall be filed in a written reasoned statement. It shall not be deemed to have been filed until the opposition fee has been paid. (Art. 99(1) European Patent Convention).

1

EP 1 599 867 B1

2

Description**Background**

[0001] Speech recognition systems, particularly computer-based speech recognition systems, are well known. Numerous inventions and voice transcription technologies have been developed to address various problems within speech recognition systems. In one aspect, advanced mathematics and processing algorithms have been developed to address the needs of translating vocal input into computer text through speech parsing, phoneme identification and database matching of the input speech so as to accurately transcribe the speech into text.

[0002] General speech recognition databases are also well known. U.S. Patent No. 6,631,348 (Wymore), for example, discloses a speech recognition system in which vocal training information is provided to create different vocal reference patterns under different ambient noise levels. The Wymore invention creates a database of captured speech from this training input. During operation, a user of the Wymore system may then dictate speech under various ambient noise conditions and the speech recognition system properly filters the noise from the user's input speech based on the different stored models to determine the appropriate, spoken words, thereby improving the accuracy of the speech transcription.

[0003] U.S. Patent No. 6,662,160 (Chien et al.) also discloses a system involving adaptive speech recognition methods that include noise compensation. Like Wymore, the system of Chien et al. neutralizes noise associated with input speech through the use of preprocessed training input. Chien et al. employs complex statistical mathematical models (e.g. Hidden Markov Models) and applies optimal equalization factors in connection with feature vectors and probability density functions related to various speech models so as to accurately recognize a user's speech.

[0004] Other voice transcription systems address the problems of minimizing and correcting misrecognition errors. For example, U.S. Patent No. 6,195,637 (Ballard et al.) discloses a transcription system that accepts a user's dictation and contemporaneously allows a user to mark misrecognized words during the dictation. At the conclusion of dictation, a computer-based, textual correction tool is invoked with which the user may correct the marked, misrecognized words. Numerous, potentially intended words, e.g. words that are close in phonetic distance to the actual speech, are provided by the Ballard et al. system for possible replacement of the misrecognized word. Other examples of misrecognized words include incorrectly spelled words and improperly formatted words, (e.g. lack of upper case, letters in a name or incorrect punctuation). In one embodiment, Ballard et al. discloses a computer having a windows-based, graphical user interface that displays the list of potentially intended words from which the user selects the appropriate word

with a graphical input device, such as a computer mouse.

[0005] Other existing speech recognition systems deal with problems associated with large, speech recognition vocabularies, i.e. the entire English language. These systems typically address the allocation of the computer-based resources required to solve the speech recognition problems associated with such a vocabulary. U.S. Patent No. 6,490,557 (Jeppesen), for example, discloses a system and method for recognizing and transcribing continuous speech in real time. In one embodiment, the disclosed speech recognition system includes multiple, geographically distributed, computer systems connected by high speed links. A portion of the disclosed computer system is responsible for preprocessing continuous speech input, such as filtering any background noise provided during the speech input, and subsequently converting the resultant speech signals into digital format. The digital signals are then transcribed into word lists upon which automatic speech recognition components operate. Jeppesen's speech recognition system is also trainable so as to accommodate more than one type of voice input, including vocal input containing different accents and dialects. Thus, this speech recognition system is capable of recognizing large vocabulary, continuous speech input in a consistent and reliable manner, particularly, speech that involves variable input rates and different dialects and accents. Jeppesen further discloses systems having on-site data storage (at the site of the speech input) and off-site data storage which stores the databases of transcribed words. Thus, in one aspect, a primary advantage of Jeppesen is that a database of large scale vocabularies containing speech dictations is distributed across different geographical areas such that users employing dialects and accents within a particular country or portion of the world would be able to use localized databases to accurately transcribe their speech input.

[0006] Other large vocabulary speech recognition systems are directed to improving the recognition of dictated input through the use of specialized, hierarchically arranged, vocabularies. The computerized, speech recognition system of U.S. Patent No. 6,528,380 (Thelan et al.), for example, employs a plurality of speech recognition models that accept incoming speech in parallel and attempts to match the speech input within specific databases. Since the English language vocabulary, for example, is relatively large, the speech matching success rate using such a large vocabulary for any given particular dictation may be lower than what is acceptable for a particular application. Thelan et al. attempts to solve this problem through the use of specific vocabularies selected by the voice recognition modules after a particular speech vocabulary and associated text database is determined to be more appropriately suited to the dictation at issue. Thus, Thelan et al. begins with an ultra-large vocabulary and narrows the text selection vocabularies depending on the speech input so as to select further refined vocabularies that provide greater transcription

3

EP 1 599 867 B1

4

accuracy. Model selectors are operative with Thelan et al. to enable the recognition of more specific models if the specific models obtain good recognition results. These specific models may then be used as replacement for the more generic vocabulary model. As with Jeppesen, Thelan et al. discloses computer-based speech recognition system having potentially distributed vocabulary databases.

[0007] Heretofore, no computerized speech recognition systems have been developed that take advantage of repeated dictation of specific terms into specific form fields or repeated dictation of specific terms by specific persons. In particular, context-specific vocabularies or context-specific modifications of matching probabilities have not been provided with respect a context specific vocabulary which is used on conjunction with more general vocabularies. The modern necessity of using specific, computerized, form-based input creates a unique problem in that the general vocabularies used by many of the commercial speech recognition software programs do not provide efficient and accurate recognition and transcription of users' input speech. The limitations of the present systems lie in the fact that any vocabulary large enough to accommodate general as well as specific text will have phonetically similar general text so as to cause an unacceptably high error rate.

[0008] WO 01/26093 describes a speech recognition system for searching a first grammar file for a matching phrase and searching a second grammar file if a matching phrase is not found in the first grammar file.

[0009] WO 01/69905 describes a speech recognition system that stores caller specific voice files to help recognise speech of frequent callers and a generic voice file used when a caller does not have a caller-specific description.

[0010] Aspects of the present invention are set out in the appended independent claims.

[0011] According to a preferred embodiment of the invention, a method for improving the accuracy of a computerized, speech recognition system, the speech recognition system including a base vocabulary, the method includes loading a specified vocabulary into computer storage, the specified vocabulary associated with a specific context; accepting a user's voice input into the speech recognition system; evaluating the user's voice input with data values from the specified vocabulary according to an evaluation criterion; selecting a particular data value as an input into a computerized form field if the evaluation criterion is met; and if the user's voice input does not meet the evaluation criterion, selecting a data value from the base vocabulary as an input into the computerized form field. According to further embodiments, the method further includes evaluating the user's voice input with data values from the base vocabulary according to a base evaluation criterion if the user's voice input does not meet the evaluation criterion. According to another embodiment, the evaluation criterion is a use weighting associated with the data values. The evaluation

ing may further include the step of applying a matching heuristic against a known threshold. The step of applying a matching heuristic may further include a step of comparing the user's voice input to a threshold probability of matching an acoustic model derived from the specified vocabulary. The context is associated with any one or more of the following: a topical subject, a specific user, and a context is associated with a field.

[0012] According to another preferred embodiment of the invention, a method for improving the accuracy of a computerized, speech recognition system is provided that include the steps of loading a first specified vocabulary into computer storage, the first specified vocabulary associated with a first computerized form field; accepting a user's voice input into the speech recognition system; evaluating the user's voice input with data values from the first specified vocabulary according to an evaluation criterion; selecting a particular data value as input into the first computerized form field if the user's voice input meets the evaluation criterion; loading a second specified vocabulary into computer storage, the second specified vocabulary associated with a second computerized form field; accepting a user's voice input into the speech recognition system; evaluating the user's voice input with against data values from the specified vocabulary according to an evaluation criterion; and selecting a particular data value as input into a second computerized form field if the user's voice input meets the evaluation criterion. In one aspect, the evaluation criterion for the steps of evaluating the first and the second specified vocabularies are the same. In another aspect, the evaluation criterion for the steps of evaluating the first and the second specified vocabularies are different criterion. In still another aspect, the first and second computerized form fields are associated with different fields of a computerized medical form.

[0013] In yet another embodiment the present invention provides a method for improving the accuracy of a computerized, speech recognition system that includes loading a first specified vocabulary into computer storage, the first specified vocabulary associated with a first user of the speech recognition system; accepting the first user's voice input into the speech recognition system; evaluating the first user's voice input with data values from the first specified vocabulary according to an evaluation criterion; selecting a particular data value as an input into a computerized form field if the first user's voice input meets the evaluation criterion; loading a second specified vocabulary into computer storage, the second specified vocabulary associated with a second user of the speech recognition system; accepting a second user's voice input into the speech recognition system; evaluating the second user's voice input with data values from the specified vocabulary according to an evaluation criterion; and selecting a particular data value as an input into the computerized form field if the second user's voice input meets the evaluation criterion. In one aspect, the first and second users of the speech recognition system

3

5

EP 1 599 867 B1

6

are different doctors and the computerized form fields are associated with a field within a computerized medical form.

[0014] In still another embodiment of the present invention, a method is provided for improving the accuracy of a computerized, speech recognition system that includes loading a first specified vocabulary into computer storage, the first specified vocabulary associated with a first context used within the speech recognition system; accepting a user's voice input into the speech recognition system; evaluating the user's voice input with data values from the first specified vocabulary according to an evaluation criterion; selecting a particular data value as an input into a computerized form field if the user's voice input meets the evaluation criterion; loading a second specified vocabulary into computer storage, the second specified vocabulary associated with a second context used within the speech recognition system; accepting the user's voice input into the speech recognition system; evaluating the user's voice input with data values from the specified vocabulary according to an evaluation criterion; and selecting a particular data value as an input into the computerized form field if the user's voice input meets the evaluation criterion. In one aspect, the first context is a patient's age and the second context is a patient diagnosis of the patient

[0015] In still another embodiment of the present invention, a computerized speech recognition system is provided including a computerized form including at least one computerized form field; a first vocabulary database containing data entries for the computerized form field, the first vocabulary associated with a specific criterion; a second vocabulary database containing data entries for the data field; and an input for accepting a user's vocal input, the vocal input being compared to the first vocabulary as a first pass in selecting an input for the computerized form field, and the vocal input being compared to the second vocabulary as a second pass in selecting an input for the computerized form field. In one aspect, the criterion is one or more of the following: a topical context, a specific user of the speech recognition system, a form field. In another aspect, the first vocabulary database is a subset of the second vocabulary database.

[0016] In yet another embodiment of the present invention, a database of data values for use in a computerized speech recognition system is provided including a first vocabulary database containing data entries for a computerized form including at least one computerized form field, the first vocabulary associated with a specific criterion; and a second vocabulary database containing data entries for the data field. In one aspect, the criterion is one or more of the following: a topical context, a specific user of the speech recognition system, a field.

Brief Description of the Drawings

[0017] The invention and its wide variety of potential embodiments will be readily understood via the following

detailed description of certain exemplary embodiments, with reference to the accompanying drawings in which:

[0018] FIG. 1 is a general network diagram of the computerized speech recognition system according to one embodiment of the present invention;

[0019] FIG. 2 is a system architecture diagram of a speech recognition system according to one embodiment of the present invention;

[0020] FIG. 3 shows an arrangement of a graphical user interface display and associated data bases according to one embodiment of the present invention;

[0021] FIG. 4 is a graphical depiction of different text string database organizations according to one embodiment of the present invention;

[0022] FIG. 5 is a graphical depiction of one specific, text string database according to one embodiment of the present invention.

[0023] FIG. 6 is a graphical depiction of another specific, text string database according to one embodiment of the present invention;

[0024] FIG. 7 is a process flow diagram for the speech recognition system according to one embodiment of the present invention; and

[0025] FIG. 8 is another process flow diagram for the speech recognition system according to another embodiment of the present invention.

Detailed Description

[0026] Specific examples of the present invention are provided within the following description. Persons of skill in the art will recognize that these are merely specific examples and that more general uses for the present invention are possible. Specifically, in the examples that follow, the present invention is generally described as it pertains to speech recognition within the medical field and as it may be used within a medical office. It is easily understood and recognized that other applications of the present invention exist in other fields of use, including use in a general web-based form, or web page. Further, the system of the present invention is described as being implemented in software, but hardware and firmware equivalents may also be realized by those skilled in the art. Finally, the pronoun, "he", will be used in the following examples to mean either "he" or "she", and "his", will be used to mean either "his" or "her".

[0027] Fig. 1 shows a general office environment including a distributed computer network for implementing the present invention according to one embodiment thereof. Medical office 100 includes computer system 105 that is running speech recognition software, microphone input 110 and associated databases and memory storage 115. The computerized system within office 1 may be used for multiple purposes within that office, one of which may be the transcription of dictation related to the use of certain medical forms within that office. Office 1 and its computer system(s) may be connected via a link 130 to the internet in general, 140. This link may

7

EP 1 599 867 B1

8

include any know or future devised connection technology including, but not limited to broadband connections, narrow band connections and/or wireless connections. Other medical offices, for example offices 2 through N, 151-153, may also be connected to one another and/or to the internet via data links 140 and thus to office 1. Each of the other medical offices may contain similar computer equipment, including computer equipment running speech recognition software, microphones, and databases. Also connected to internet 140 via data link 162 is data storage facility 170 containing one or more speech recognition databases for use with the present invention.

[0028] Fig. 2 provides a diagram of a high-level system architecture for the speech recognition system 200 according to one embodiment of the present invention. It should be recognized that any one of the individual pieces and/or subsets of the system architecture may be distributed and contained within any one or more of the various offices or data storage facilities provided in Fig. 1. Thus, there is no preconceived restriction on where any one of the individual components within Fig. 2 resides, and those of skill in the art will recognize various advantages by including the particular components provided in Fig. 2 in particular geographic and data-centric locations shown in Fig. 1.

[0029] Referring to Fig. 2, input speech 205 is provided to the speech recognition system via a voice collection device, for example, a microphone 210. Microphone 210 in turn is connected to the computer equipment associated with the microphone, shown as 105 in Fig. 1. Computer system 105 also includes a speech recognition software system 212. Numerous, commercial speech recognition software systems are readily available for such purpose including, but not limited to, ViaVoice offered by IBM and Dragon Naturally Speaking offered by ScanSoft. Regardless of the manufacturer of the product, the speech recognition software includes, generally, a speech recognition module 217 which is responsible for parsing the input speech 205 as digitized by the microphone 210 according to various, well-known speech recognition algorithms and heuristics. Language model 219 is also typically included with speech recognition software 212. In part, the language model 219 is responsible for parsing the input speech according to various algorithms and producing fundamental language components. These language components are typically created in relation to a particular language and/or application of interest, which the speech recognition system then evaluates against a textual vocabulary database 220 to determine a match. In frame-based systems, for example, incoming analog speech is digitized and the amplitude of different frequency bands are stored as dimensions of a vector. This is performed for each of between 6,000 and 16,000 frames per second and the resulting temporal sequence of vectors is converted, by any of various means, to a series of temporally overlapping "tokens" as defined in U. S. Pat. # 6,073,097.

These tokens are then matches with similar temporal se-

quences of vectors generated from strings of text in the active vocabulary according to the active language model and any active set of "learned" user-specific phonetic patterns and habits.

[0030] General text database 220 is typically included as part of speech recognition software 212 and includes language text that is output by the speech recognition software once a match with the input speech is made. General or base vocabulary database 220 may contain the textual vocabulary for an entire language, e.g. English. More typical, however, the base vocabulary database contains a sizable subset of a particular language or desired application, e.g. hundreds of thousands of words. Those of skill in the arts of database management and computer science will realize that certain inherent computational difficulties and computer processing problems exist in the use and management of databases of this size. The principal barrier to accurate speech matching (recognition) with large vocabularies is "background noise" in the form of sufficient numbers of phonetically similar text mismatches in the vocabulary to give an unacceptable frequency of transcription errors. Other problems include the latency associated with full database searches for textual matches corresponding to input speech and the time and computer processing resources that must be expended within applications in which the base vocabulary database is swappable and must be replaced. These problems will arise, for example, with rapid swapping of large vocabulary databases in different languages.

[0031] Following a textual match from the speech input by speech recognition system 212, the text output from base vocabulary database 220 is then provided as input to any one of a number of other computer-based applications 230 into which the user desires the text. Examples of typical computer applications that are particularly suited for use with speech recognition software include, but are not limited to word processors, spreadsheets, command systems and/or transcription systems that can take advantage of a user's vocal input. Alternatively, as more text-based applications accompany people's use of the internet, for example, such vocal input may be used to provide inputs to text field within a particular form, field or web page displayed by an internet browser.

[0032] Although the initial applications of the present invention are directed to voice-to-text applications in which vocal input is provided and textual output is desired, other applications are envisioned in which any user or machine provides an input to a recognition system, and that recognition system provides some type of output from a library of possible outputs. Examples of such applications include, but are not limited to a search and match of graphical outputs based on a user's voice input or an action-based output (e.g. a computer logon) based on a vocal input. One example of an action-based output may be to provide access to one of several computer systems, the list of computer systems being stored in a database of all accessible computer systems based on

a user's bio-input (e.g. fingerprint) or a machine's mechanical input (e.g. a login message from a computer).

[0033] Referring again to Fig. 2, the speech recognition/voice transcription system further includes a specified database of text string values that provide a first-pass output in response to a particular speech input against which the system attempts to determine a match. These text strings may be stored in any one of a number of formats and may be organized in any one of a number of manners depending on the practical application of the system. In one particularly preferred embodiment, the text strings within specified database 250 are provided from the vocal inputs of previous users of the speech recognition system. Using the Doctor's office example shown in of Fig. 1, the first-pass text strings may be organized by users (e.g. doctors) of the system such that those text strings used by a particular doctor are loaded by the system as first-pass potential matches when that particular doctor logs into the system and/or his vocal speech is recognized and identified by the system as belonging to that doctor. Sub-databases 261, 262 and 263 illustrate such an organization based on users of the system.

[0034] Specified database 250 may also be organized according to numerous other criteria that may be advantageous to users of the speech recognition system. In another arrangement, the sub-databases of first-pass text strings within first-pass, specified database 250 may be organized by fields within a computerized or web-based electronic form. Using the example of a doctor's office once again and referring to Fig. 3, text input may need to be input into a medical form 310, that includes a patient's name, shown in computerized form field 315, the patient's address, shown in computerized form field 318, the patient's phone number, shown in computerized form field 320, and the patient's age, shown in computerized form field 320. Sub-databases 371, 372 and 373 shown in Fig. 3 are specific examples of the general field sub-databases 271, 272 and 273 of Fig. 2. These sub-databases provide first-pass text strings for matching speech input provided by the doctor when populating form fields 315, 318 and 328 (Fig. 3) respectively.

[0035] As yet another example of sub-database organization within specified database 250, a context associated with some aspect of the present speech input (or even past speech input) may be used to organize and condition the data into appropriate first-pass sub-databases. For example, the sub-database 381 associated with the findings field 330 within the medical form of Fig. 3 may be conditioned upon both the history and the age of the patient under the presumption that previous findings related to a particular combination of history and age group, either within an individual medical office or in general, are more likely to be repeated in future speech inputs with respect to patients having the same combination of age range and history. As one example, the findings fields populated within a form in the office practice of a primary care physician, with a history of abdominal pain and char-

acteristic physical findings may be quite similar for the following two conditions: "appendicitis" as a probable "Interpretation" field for patients age 5-12; and "diverticulitis" as a probable "Interpretation" for patients age 75+. Characteristic findings (abdominal pain with what is called "rebound tenderness") will be stored in sub-database 381 and provided to "findings" field 330, while "appendicitis" and "diverticulitis" will be stored in sub-database 382 and provided to "Interpretation" field 350.

[0036] Specified database 250 may be created and organized in any number of ways and from any one of a number of sources of information so as to provide an accurate first-pass database for appropriate and efficient use within a particular context. If, for example, specified database 250 contains text strings organized by users of the system (a user context) under the statistical presumption that each specific doctor is more likely to repeat his or her own relatively recent utterances than earlier utterances, in situations when all other system parameters are the same, and more likely to repeat terms used by other system users or other physicians in the same specialty under otherwise identical circumstances, than to use terms neither they nor others have used in that situation, text from their own past dictations or those of others (whether manually or electronically transcribed) may be used to populate and arrange the text string values within the database. If, however, a high probability first-pass database is used to provide text strings to be input into particular fields within a computerized form, then these data values may be derived and input from previously filled-out forms. These data may then be organized into sub-databases according to form fields, for example as shown in Fig. 3 by sub-databases 371-381. Also, the specified database 250 may contain one, many or all such data for use within a particular desired context and output application. Finally, the actual data values within the database may be dynamically updated and rearranged into different sub-databases during the actual use of the speech recognition system so as to accommodate any particularly desirable speech recognition situation. In the most useful instances, the data values that populate the specified database 250 will be obtained from historical data and text strings that accompany a particular use and application of the speech recognition system.

[0037] Supplemental data may also accompany the data values and text strings stored within specified database 250. In particular, weightings and prioritization information may be included as part of the textual data records that are to be matched to the input speech. These weightings may help determine which data values are selected, when several possible data values are matched as possible outputs in response to a particular speech input. Further, these weighting and prioritization information may be dynamically updated during the course of the operation of the speech recognition system to reflect prior speech input. Those of skill in the art will realize a plurality of ways in which the data elements within the

11

EP 1 599 867 B1

12

specified database may be rearranged and conditioned so as to provide an optimal first-pass database for use in the speech recognition system of the present invention.

[0038] Referring again to Fig. 2, the speech recognition/voice transcription system of the present invention further includes a context identification module 240. The context identification module is coupled to one or more input and recognition components (Fig. 2, 205-230) of the overall speech recognition system 200 and is used to select or create a proper sub-database within the entire specified database 250. If, for example, the desired sub-databases to be used are based on a user context, then the context identification module may take input from a user identification device (not shown) or may determine the user from speech characteristics determined by the speech recognition software so as to select an appropriate user sub-database (e.g. 261) from the entire specified database 250. Alternatively, the data values within the specified database 250 may be loosely organized and the context identification module may actually condition the data values so as to dynamically create an appropriate user sub-database from the information stored within the specified database. As another example, the context identification module may monitor and interpret a particular form field that is active within an application 230 into which text input is to be provided. After making such a determination, the context identification module may select, or as mentioned above, actually condition the data values so as to dynamically create, an appropriate user sub-database from the information stored within the specified database.

[0039] Referring again to Fig. 2, the speech recognition/voice transcription system may further include a prioritization module 245. As with the context identification module, the prioritization module may be coupled to any one or more input and recognition components (Fig. 2, 205-230) within the overall speech recognition system 200 including the specified database 250. As mentioned above and provided in more detail below, the prioritization module assists in collecting actual use information from the speech recognition system and using that data to dynamically prioritize the data values within any or all of the sub-databases contained within specified database 250.

[0040] In one particularly preferred embodiment of the present invention, specified database 250 contains text strings as selectable data values for input into medical forms in a word processing application 230. The text strings may be organized according to a number of different criteria based on the users of the forms and/or the fields within the electronic forms. As shown in Fig. 3, a computer-based electronic medical form 310 shows several fields within a medical report. For example, computerized electronic form 310 may include a name field 315, an address field 318, a phone number field 320, as well as more general fields such as a findings field 330 and an interpretations field 350. One possible organization of the text string data values within specified database

250 is to associate each text string with each field within a particular electronic form. As shown in Fig. 3, text string sub-database 371 may be associated with name field 315, text string sub-database 372 may be associated with address field 318 and text string sub-database 381 may be associated with findings field 330. In this particular example, two separate organizations of the text strings exist within specified, text string sub-databases 371 through 382. For single context fields, the name field 315 for example, sub-database 371 may contain text strings that only indicate patient's names. Likewise, text string sub-database 372 associated with address field 318 of electronic computer form 310 may contain only text strings associated with street addresses.

[0041] It should be noted that the data organizations referenced by 261-283 in Fig. 2 and 379-382 in Fig. 3 are logical organizations only. The data records within specified database 250 may be organized, arranged and interrelated in any one of a number of ways, two of which are shown in Fig. 4. Referring to Fig. 4, the organization of the records within specified database 450 may be loose, i.e. all records may be within one file 455 where each record (and output text string) contains a plethora of relational information. (Option A). The relational information within the singular file would then, presumably, be able to be used to create the logical divisions shown in Figs. 2 and 3. One example of a sub-database might be a field context sub-database 471, for example, where the relational data pertaining to the form field within file 455 is used to organize the sub-database. Alternatively, organization of the records within specified database 250 may be tight, i.e. records (and output text strings) may be highly organized according to context/field/user such that a one-to-one relationship exists between a particular file of records (sub-database) and a form field or user, as shown in option B of Fig. 4. While the organization provided in option B may require more computer memory because of the information redundancy needed to create all the discrete sub-databases, this disadvantage in the overall database size 450 may be offset by the advantage of having smaller physical files 456-458 that can be more quickly swapped in and out of computer memory within the speech recognition system. In general, those of skill in the art will realize that different organizations of the same data will provide various advantages and that such data may be organized to optimize any one of number of parameters and/or the overall system operation so as to enhance the advantages of the present embodiment. Finally, a combination of both database organizations could be used to provide a system that has the advantages of the present embodiment.

[0042] Regardless of the data organization of specified database 250, two types of specified, sub-databases are contemplated. The first type may be classified as a singular context sub-database in that one specific criterion provides the motivation for grouping and organizing the records to create the sub-database. One specific embodiment of the specified, this type of sub-database, 371 of

13

EP 1 599 867 B1

14

Fig. 3, is shown in more detail in Fig. 5, where text string records containing street addresses are stored within sub-database 571 in tabular format. In this particular embodiment, individual records 510, 511 and 512 contain text strings of previously dictated (specified) street addresses which are provided for the purpose of matching a user's speech input when the address field 318 (Fig.3) is the active dictation field. Other data, such as weighting information 552 and user's data 554, may also be included within text string sub-database 371. With reference to the specific example of Fig. 5, the data records within the sub-database 571 contain text strings and accompanying relational data intended for use only within a specific field within a computerized form or web page. Other specified sub-databases similar to 571 may contain text strings and accompanying relational data that is intended for use with only one of the users of the speech recognition system.

[0043] In a second sub-database type, multiple context organizations of the data within specified database 250 are also created. For example, medical form 310 of Fig. 3 may contain input fields that are related to other input fields within the overall electronic form. This interrelationship typically occurs when the voice dictation provided as an input to a field within an electronic form is of a more general nature. In particular, the organization of the text strings within a sub-database may not be based on a single, external, context, such as a specific user of the system or a particular field within an electronic form, but rather may be based on the interrelation of the actual text strings in a more complex manner. As one example, context specific sub-databases 381 (pertaining to the medical findings field) and 382 (pertaining to the medical interpretations field) may include contextually intertwined text strings that the speech recognition system of the present invention must identify and properly select so as to achieve the efficiencies of the present embodiment. These more complex, contextually, intertwined text string sub-databases are shown as logical sub-databases 281-283 in Fig. 2.

[0044] A simplified example of the above-mentioned text string interrelation is provided below. As shown in Fig. 3, sub-database 381 provides text strings that may be input into findings field 330 and sub-database 382 provides text strings that may be input into interpretations field 350. However, unlike fields with a limited range of accepted input within the electronic computer form, the name field 315 for example, sub-database 381 is designed to match text strings to a more general and varied voice input provided to the speech recognition system. Fig. 6 shows one specific embodiment of the specified, text string sub-database 382 of Fig. 3. Sub-database 382 provides text string records related to medical interpretations which are stored within sub-database 682 in tabular format. In this particular embodiment, individual records 615, 616 and 617 contain text strings from previously dictated (specified) interpretations which are provided for the purpose of matching a user's speech input

when the interpretations field 350 (Fig.3) is the active dictation field. Other relational data, such as weighting information 652 and interrelational context information (e.g. age 654, user 656, findings 658) may also be included within text string sub-database 682. In the example of Fig. 6, interpretations text strings, such as pneumonia and dysphagia, are provided as potential text strings to be evaluated against a user's dictation to provide a text input to the interpretations field.

[0045] Also shown in Fig. 6 are, two, similar sounding medical terms that have entirely different meanings: dysphagia - a difficulty in swallowing, and dysphasia - an impairment of speech consisting in lack of coordination and failure to arrange words in a proper order. The interpretations sub-database 682 includes both textual inputs as records 616 and 617 respectively. Exemplary interrelational data are also included as data within the text records record of the sub-database. Such data include a patient's history 654, a user of the system 658, the specific findings regarding the patient 658, as well as a general, historical weighting based on the number of times the two term have been used 652. During a dictation into the interpretations field 350 of electronic form 310, table 682 is loaded and consulted to achieve the best possible textual input for dictated speech. If, for example, the phonetically similar word dysphagia/dysphasia is dictated into the system then the context interpretation module would evaluate that voice input in view of any one or combination of contextual data. In one case, if the patient's past medical history included digestive complaints then the more probable textual match, dysphagia, may be selected. Similarly, if the patient's past medical history included neurological complaints, the term dysphasia may be selected. Similarly, the context identification module may rely upon other relational data associated with the two text strings to determine the highest probability input. If Dr. Brown is a pediatrician and Dr. Smith is a geriatric physician, then appropriate weight may also be given by the selection system to these previous inputs in determining the proper text input for the interpretations field. Likewise, the input to the findings field 330 may be considered, in which a "difficulty swallowing" would result in a more likely match with dysphagia and "speech impairment" would result in a more likely indication of dysphasia. In addition, other simple weighting factors such as the number of times each term has been used previously may also be used by the system of the present invention to select a more probable input text string. Finally, the system may use one, many, or all of the aforementioned contextual relationships to determine and select the proper text input, possibly after assigning additional weighting function to the interrelational data itself, i.e. weighting a user's context higher than the age context.

[0046] In operation, a user of the speech recognition system of the embodiment inputs speech 205 to microphone 210 for processing by speech recognition system 212. As a stand-alone system, speech recognition sys-

15

EP 1 599 867 B1

16

term package 212 typically provides a single, general or base vocabulary database 220 that acts as a first and only database. Because of the size of the database and the general nature of the language and the text strings contained within it, voice-to-text transcription accuracies may vary when the speech recognition system is used only with such large, non-specific vocabularies. In medical contexts, for example, inaccuracies in transcription of dictation may result in undesirable or even disastrous consequences. Thus, the inaccuracies generally tolerated by system users must be improved. Greater transcription accuracy, as well as consistency in the dictation within fields of an electronic, computer-based form, for example, may be achieved through the use of multiple databases containing text strings previously used in different contexts. Specifically, through the proper selection of a first-pass database containing a limited but specialized vocabulary and the insertion of this first-pass database into the existing processing used by commercial voice transcription systems, the transcription accuracies of these systems can be markedly improved. Failing a match in the more specific, first-pass database, the speech recognition system can always default to the more general, base vocabulary to provide a textual match for the dictated input.

[0047] According to various embodiments of the present invention, the specified database 250 is used by the speech recognition system as a first-pass database in selecting an appropriate textual match to the input speech 205. The context identification module 240 is responsible for selecting and loading (or creating) a particular sub-database from specified database 250 during a user's dictation so as to provide a high probability of a "hit" within that sub-database. The selection process employed by context identification module is based on a context of the input speech or a context within the dictation environment. Possible contexts include, but are not limited to, a particular user of the speech recognition system, a particular field within an electronic form being processed by the speech recognition system, or the interrelation of previously input text with a sub-database of text that is likely to be dictated based thereon.

[0048] Thus, the inherent value of specified database 250 lies in its historical precedent as optionally conditioned with weighting functions that are applied to the text strings within the database. Thus, the creation of a specified database is central to its effective use within the speech recognition system.

[0049] Specified database 250 may be created in any of a number of manners. In one particularly preferred embodiment, past forms may be scanned and digitally input into a computer system such that all the text strings used within those computer forms are digitized, parsed and then stored within the database. The text strings may then be subdivided into specific databases that are applicable to specific speech recognition circumstances. For example, with respect to the example of addresses sub-database shown in Fig. 5, a series of previously re-

corded paper or electronic medical forms may be parsed, separated and stored such that all the street addresses used on those forms are stored in a separate portion 271 of database 250. Likewise, findings within field 330 and interpretations within field 330 of the electronic form in Fig. 3 may be subdivided from general text string database 250 to create a specific contextual database of diagnoses for use with a particular medical form. As previously described, those of skill in the art will recognize that specified database 250 may be organized in any one of a number of different ways to suit the particular needs of a particular speech recognition application, such as textual input into an electronic form. Such organization may take place statically, i.e. before the user employs the voice transcription system, or dynamically, i.e. during the use of the voice transcription system. In the dynamic context, certain relationships among sub-databases may also be leveraged to provide inputs between computerized form fields.

[0050] Referring to Fig. 7, a general process flow is provided for the operation of speech recognition system 200. The process starts with step 705 in which the speech recognition system is loaded and has begun to operate. Specified vocabulary databases may be defined and loaded here for a particular, more global use during the remainder of this process. Next, a user of the system is identified at step 707. As one example, the user may be a particular doctor who wishes to provide speech input to a medical form as part of his practice within a practice group or a medical office. As described above, this user ID may later be used to select appropriate sub-databases and associated text strings from specified database 250. User identification may be done through speech recognition, keyboard entry, fingerprinting or by any means presently known or heretofore developed. Next, voice input from the user is provided to the speech recognition system in step 710. This vocal input is digitized for use within computer system 105 which is then input into the speech recognition system employed on that computer system as shown in step 720.

[0051] Next, the context identification module selects or creates an appropriate sub-database consisting of a subset of the text strings within database 250 as the system's operative first-pass database at step 730. As described above, the selection of an appropriate sub-database may occur according to any one or more of a number of different criteria. In one particularly preferred embodiment, the criterion on which the sub-database is selected is based upon the user of the voice transcription system as provided in step 707. Specifically, any particular user may have a historical use of certain words and phrases which may serve as a higher probability first-pass source of text string data for future use by that particular user. Thus, the appropriate selection of that database will result in higher transcription accuracy and use within the speech recognition system.

[0052] According to another particularly preferred embodiment of the present invention, the sub-database is

17

EP 1 599 867 B1

18

selected from the specified database 250 at step 730 according to the field within the electronic form into which text is being input. For example, referring to Fig. 3, when a user wishes to populate address field 318 with a particular address, the user would indicate to the system at step 730 (e.g. through a computer graphical user interface or a vocal command input) that the address field is to be populated. The speech recognition software then selects or creates an appropriate sub-database from specified database 250 that contains at least the addresses for use within that form field. The actual data selected and pulled by the context identification module, as mentioned above, would typically include related contextual information that would provide insight into the historical use of particular addresses so as to provide a higher probability in transcription accuracy.

[0053] Referring back to Fig. 7, the speech input provided by the user to the speech recognition system at step 720 is evaluated by that system with respect to the text strings within the sub-database selected in step 740. This evaluation may be performed according to the same algorithms and processes used within the speech recognition system 212 which are used to select matching text from its own base vocabulary in database 220. Various methods and mechanisms by which the input speech is parsed and converted to a language output and/or text string output are well-known in the art, and these text matching mechanisms and evaluation criteria are independent of the other aspects of the present invention. Furthermore, other known evaluation criteria may be used on the overall database 250 or the sub-database selected in step 730. Such evaluation methods are well-known, although particular evaluation criteria that are applicable to speech recognition principles may also be employed when populating a field within an electronic form. As an example, the specific text strings of a particular sub-database, such as that shown in Fig. 5 may include a weighting function as shown in field 552 of sub-database 571. The weighting field, for example, may include the number of times a particular address has been input into a form within a specific historical period. Even with this over-simplified weighting scheme, ambiguities as between two very similar addresses may be easily resolved in determining a proper textual match corresponding to a speech input. Other weighting schemes, using both objective indicia (e.g. data use count) and subjective indicia (e.g. weights related to the data itself and its interrelation with other data) are well known in the art and may also be included within specific database 571 for use in the context identification module. Further, other evaluation criteria may be used to select an input text string from the sub-database. For example, a most-recently-used algorithm may be used to select data that may be more pertinent with respect to a particular transcription. Other weighting and evaluation criteria are well-known and those of skill in the art will appreciate different ways to organize and prioritize the data so as to achieve optimal transcription accuracy. Finally, a prioritization

module 245 may be included as part of the speech recognition system 200 to implement and manage the above-mentioned weighting and prioritization functions.

[0054] If the evaluation of the voice input at step 740 results in a match within the selected sub-database of text strings according to the evaluation criterion, then that text string is selected as an output at step 750 and the text string is used to populate the desired field within the electronic form at step 760. Alternatively, if the evaluation criteria is not met at step 740, the speech recognition system would default to base vocabulary database 220 at step 770, at which point, the speech recognition software would transcribe the user's voice input in its usual fashion to select a text string output (step 750) according to its own best recognition principles and output the same to the electronic form (step 760).

[0055] It should be recognized that the steps provided in Fig. 7 may be repetitively performed in a number of different ways. For example, as one particular electronic form is filled out, sequential fields within that form need to be designated and then populated with an appropriate text string. As such, following the insertion of a particular text string within a particular form field, the process of Fig. 7 may return to step 720 where the user inputs additional speech input after selecting the new field into which the vocal input is to be transcribed. During this second iteration, a second, appropriate sub-database of text strings from specified database 250 would be selected as an appropriate first-pass database for the second field. The process of evaluating and matching the user's vocal input with text strings within the second sub-database, i.e., steps 740 through 770, would operate as mentioned above.

[0056] In another operative alternative, a second user may employ the speech recognition system in response to which different sub-databases of text strings would be loaded that pertain to the specific use of that second user at step 730. In this iterative process, a second user would be identified at step 707, after which the speech input provided by that second user would be digitized and processed by the speech recognition system at step 720. The selection and/or creation step 730 may or may not be performed (again) and may be omitted if the only sub-database selection step is conditioned upon a user. The remainder of the process provided in Fig. 7 may then be performed to select an appropriate text string as input into the fields of the electronic form for that second user.

[0057] Specific scenarios in which the embodiments might be used in a medical office are provided below.

[0058] Example #1: A new radiologist joins a group of radiologists who have been using voice recognition technology to dictate reports for about two years. Their practice has a four year old database of digitally recorded imaging studies, linked to a database of the past two years of computer-transcribed reports as well as several years of prior reports manually transcribed to computer by transcriptionists listening to voice recordings. The new radiologist has "trained" the voice engine to recognize

19

EP 1 599 867 B1

20

his voice as a new user by engaging in a set of radiology voice training exercises that are customized to include phrases commonly used by other members of his group.

[0059] If the new radiologist's first assignment using the system is to dictate a report on a sinus CT scan, the radiologist would identify this report as being for a sinus CT scan and click on the "findings" field at which time the program will load a specified vocabulary for first pass pre-screening composed of text strings that other members of the group have previously used in their dictations as input to the "findings" field for sinus CT scans.

[0060] Since the new radiologist is more likely to use terms previously used by his colleagues in dictating reports of previous sinus CT scans than other x-ray related terms that are phonetically similar, pre-screening the new radiologist's dictation to match text strings previously used by his colleagues, for example, in the "findings" field, will deliver a higher transcription accuracy than the use of a general radiology dictionary or a full English language vocabulary. This is so even if the general radiology vocabulary has been enriched by "learning" the preferred terminology and syntax of his colleagues. When the radiologist advances to the "interpretations" field, the virtual vocabulary previously loaded for the "findings" field will be unloaded and replaced by a similarly selected virtual vocabulary for the "interpretations" field.

[0061] As the new radiologist uses the system, the prioritization algorithm administered by the prioritization module for his specific user sub-database files may assign relatively higher prioritization scores to his own dictated text strings vis-a-vis the dictated text of his colleagues. Over time it will adapt to his personal style, further improving transcription accuracy.

[0062] Assume that on his second day of work, the new radiologist is assigned to read studies of the digestive system, and his first two cases are barium swallow studies of the upper gastrointestinal tract. The first case is for the evaluation of a two-month old infant suffering from vomiting, and the second case is a follow-up study for an 87 year-old man with esophageal strictures. While the study is the same, his findings and interpretations in the two cases are likely to be different. Depending on the number of prior reports in his practice group's database, the transcription accuracy of the new radiologist's reports may be maximized by applying more complex prioritization and selection algorithms to the selection of previously-used phrases to be loaded for first pass pre-screening. The weighting of previously used text strings and the selection of those data items as first-pass text strings values for these reports could result in the assignment of multipliers to those data items. These weights could be updated not only each time the first-pass text strings were previously used but also based on the type of study, the age of patient and the diagnoses or symptoms listed as reasons for physician's request in ordered the study. For the above-mentioned infant, weighting factors for text string prioritization and selection could, for example, be based on prior frequency of use in reports of all barium

swallow studies in children aged less than 6 months or less than one year. For the 87 year old man, such prioritization could, for example, be based on the frequency of use of those text strings in reporting barium swallow studies in patients in any one or more of the following classes: patients more than age 60/70/80; use of those text strings in reporting barium swallow studies in males in these age ranges; prior use of those text strings in reporting barium swallow studies in patients with a prior diagnosis of esophageal stricture; prior use of those text strings in reporting barium swallow studies of patients with a prior diagnosis of esophageal stricture by age and/or sex; and/or the presence or absence of other symptoms (such as swallowing pain or vomiting). Finally, the weighting factors related to the presence or absence of a symptom, including associated diagnoses (such as status post radiation therapy for a specific type of lung cancer) may be listed in the ordering physician's request for the procedure or may already be present in the database of prior diagnoses for that patient.

[0063] There may be an increased likelihood that text strings will be used in a radiology report if they have previously been used in reporting the same type of study or a related study for the same patient (as when high resolution chest tomography is ordered as a follow up to an abnormal chest x-ray). Dictation transcription accuracy may thus be improved by a prioritization algorithm that assigns increased weight to text strings that are previously used in reporting studies with these types of relationship to a study currently being conducted.

[0064] The larger the group of users that share common data and voice match text string sources, the greater the extent to which increasingly complex prioritization algorithms can increase transcription accuracy. In certain context driven applications, such as dictations related to the practice of medicine, the greater the linkage of source dictated text to the text strings from which it came, the better the ability to retrospectively analyze prioritization algorithm performance and compare the efficiency of the first-pass vocabulary based on different weighting assignments for different factors in the prioritization algorithm. This makes it possible to create first-pass databases for user in large installations, as they accumulate data with use, thereby allowing complex prioritization algorithms, to be optimized based on their own prior experiences.

[0065] Example #2: A physician dictates into either a computerized medical record database or a structured consultation report form as he examines a patient in an office setting. In this scenario, the medical report will usually begin With a listing of the problem(s) for which patient is being seen. These factors, in addition to age and sex, serve as effective weighting factors so as to allow the prioritization of previously-used text strings and load the most probable first-pass text strings for each report. Previous diagnoses, if noted in an initial consultation or if already present in the database from previous diagnosis of the same patient, may also be useful as text string

21

EP 1 599 867 B1

22

weighting factors for sub-database prioritization and selection. If the patient has been previously seen and his or her own previous reports are included in the same database, it may be efficient to assign a first multiplier or weighting factor to every prior text string used in previous reports for that patient and another multiplier or weighting factor for each text string uses in the reports for which each specific diagnosis is listed among the reasons or problems assessed at this visit.

[0066] With respect to electronic forms, a computerized medical record has functionally separate data fields. In addition, other types of medical reports have structured sections. Speech recognition transcription accuracy for each such application can be enhanced through the prioritization and selection of first pass, text string databases for each such field on the basis of numerous factors including, but not limited to: the age and sex of the patient; problems listed as reason for that patient's visit or to be determined during that patient's visit; previously recorded diagnoses for that patient; previous use of text strings to be prioritized by that physician in reports for that patient; previous use of those text strings with that combination of other selection factors by that physician for other patients; and/or previous use with that combination of other factors by other members of that specialty.

[0067] As in Example #1, as each office that uses the embodiment accumulates data, it becomes possible to retrospectively analyze prioritization algorithm performance and compare the first-pass hit efficiency of different weighting assignments for different factors in the prioritization algorithm. This allows the initial data record selection scheme to be optimized and permits for a quantitative analysis of the relative efficiency of various prioritization models and weightings for the various offices.

[0068] The specific embodiment of the present invention provided above is somewhat idealistic in that it presumes that commercially available speech recognition software provides for dynamically loadable databases and the possibility to hierarchically direct the speech recognition software to sequentially search several such loaded databases, including possibly the general or base vocabulary that the software is programmed to operate with for most other dictations. Unfortunately, none of the speech recognition software packages examined include these general capabilities. Thus, certain improvisations have been made with respect to an existing speech recognition software package in order to practice the advantages of the embodiment as described below.

[0069] In one particular application, the speech recognition software interfaces with computer operating systems according to an industry standard called the "Speech Application Programming Interface" protocol, abbreviated, "SAPI." SAPI was originally designed for the Microsoft™ Windows operating systems. During the 1990's a similar protocol called SRAPI was developed for non-Windows operating systems but SRAPI lost support in the computer industry and current versions of SAPI have been applied to non-Windows as well as Windows

operating systems.

[0070] SAPI (and, in its day, SRAPI) provide for computer-based responses to three types of speech input application defined commands, user-defined commands (both referred to hereinafter as "commands") and general dictation of vocabulary. A signal representing an incoming item of speech is first screened by the program to see if it represents a command, such as, "New paragraph," and, if so, executes it as such. Within speech recognition applications such as a word processor, this command may cause the insertion of a paragraph break, a new-line feed and an indent so as to permit the continued dictation in a new paragraph. Incoming speech items that are not recognized as commands are transcribed as general vocabulary text, in which the speech recognition software looks for the best possible match for the dictated text within combinations of single word text strings loaded into the general vocabulary database of the application.

[0071] Current versions of the SAPI protocol and current voice engines only accommodate the loading of one vocabulary at a time. However, they accept rapid loading and unloading of smaller sets of user-defined commands. These smaller sets may be as large as the relatively small, first-pass vocabularies needed to optimize speech recognition accuracy for dictation into a computer field. The embodiment comprises methods to identify, prioritize and select the high probability text strings which would optimize transcription accuracy if used as a first pass pre-screening vocabulary. These text strings may then be translated into user-defined commands which are loaded and screened for matches as a first pass "virtual vocabulary." In this manner, the existing speech recognition systems have been tricked into implementing a two-pass vocabulary screening model as described above under present SAPI protocols and with presently available voice engines. Incorporation of the methods and apparatus of the embodiment; would be made more user-friendly by incorporating the entirety of this embodiment into future versions of SAPI and into applications compliant with such future versions of SAPI.

[0072] Referring to Fig. 8, a general process flow for the operation of the speech recognition system 200 is provided as it would be implemented within a specific SAPI speech recognition engine. In general, the steps are substantially similar to those provided in Fig. 7 with the following modifications. At step 740, instead of evaluating the speech input against a set of text strings in the selected/created database, the process of Fig. 8 sequentially evaluates the speech input first, against the database of system commands 840, and then, if necessary, against the database of user-defined commands 841, and then, if necessary, against the database of a first vocabulary 842, and then, if necessary, against the database of a second vocabulary 842, and finally, if necessary, against a final database 844. If a match is determined during any one of these evaluations (steps 850-853), then either the "command" is executed (steps 854-855) or a learning function is exercised (steps

23

EP 1 599 867 B1

24

856-858), and the executed command or selected text from a database results in the generation and insertion of the selected text string into a computer form field (step 860).

[0073] With specific application to Example #1 provided above, the method provided in the flow diagram of Fig. 7 may be modified to operation more efficiently by including some of the elements of the process shown in Fig. 8. For each context of user (radiologist), type of imaging study (as chest x-ray or sinus CT), patient demographics (including age, sex, past medical history, reason for this study) and field of report, first pass vocabulary 842 may be provided which includes previous dictations by the same user when all the other variables were identical. The second pass vocabulary 843 may be provided which includes dictations by other members of the radiology group when all other variables were the same as those of the present report. The third pass vocabulary 844 may be provided which includes other dictations by the present radiologist into the same field for the same type of study but for patients with all combinations of age, sex, past medical history and reason for study. Thus a multiple pass series of specific context dependant sub-databases may be provided in actual application before the base vocabulary of the speech recognition software is employed to provide a match.

[0074] Although the invention has been described with reference to specific exemplary embodiments thereof, it will be understood that numerous variations, modifications and additional embodiments are possible, and accordingly, all such variations, modifications, and embodiments are to be regarded as being within the scope of the invention. As such, the intended scope of the invention is intended to be limited only by the claims of the invention and not by any one aspect of the description provided above since the drawings and descriptions are to be regarded as illustrative in nature only.

Claims

1. A method of operating a speech recognition system, said speech recognition system including a base vocabulary, the method comprising:

creating a specified database containing text strings, wherein the text strings are provided from the inputs of previous use of the system;
 defining a sub-database within the specified database containing text strings associated with a context of input data;
 identifying the context of an input of data;
 loading (705) a specified vocabulary from the sub-database into computer storage, said specified vocabulary associated with a specific context;
 accepting (710) a user's voice input into said speech recognition system;

evaluating (740) said user's voice input with data values from said specified vocabulary according to an evaluation criterion;

selecting a particular data value as an input into a computerised form field if said evaluation criterion is met; and
 if said user's voice input does not meet said evaluation criterion, selecting a data value from said base vocabulary as an input into said computerised form field.

2. The speech recognition method of claim 1 wherein said context is associated with said field.

3. The speech recognition method of claim 1 wherein said context is associated with a topical subject.

4. The speech recognition method of claim 1 wherein said context is associated with a specific user.

5. The method of claim 1 further comprising evaluating said user's voice input with data values from said base vocabulary according to a base evaluation criterion if said user's voice input does not meet said evaluation criterion.

6. The method of claim 1 wherein said evaluation criterion is a use weighting associated with said data values.

7. The method of claim 1 wherein said step of evaluating further includes the step of applying a matching heuristic against a known threshold.

8. The method of claim 7 wherein said step of applying a matching heuristic further includes a step of comparing said user's voice input to a threshold probability of matching an acoustic model derived from said specified vocabulary.

9. A method as claimed in claim 1, said first specified vocabulary is associated with a first computerised form field;

loading a second specified vocabulary from a second sub-database into computer storage, said second specified vocabulary associated with a second computerised form field;

accepting a user's further voice input into said speech recognition system;

evaluating said user's voice input with data values from said specified vocabulary according to an evaluation criterion; and

selecting a particular data value as input into a second computerised form field if said user's voice input meets said evaluation criterion.

10. The method of claim 9 wherein said evaluation criterion for said steps of evaluating said first and said

25

EP 1 599 867 B1

26

second specified vocabularies are the same.

11. The method of claim 9 wherein said evaluation criterion for said steps of evaluating said first and said second specified vocabularies are different criterion. 5

12. The method of claim 9 wherein said first and second computerised form fields are associated with different fields of a computerised medical form. 10

13. A method as claimed in claim 1, comprising:

loading a second specified vocabulary from a second sub-database into computer storage, said second specified vocabulary associated with a second user of said speech recognition system ; 15
accepting a second user's voice input into said speech recognition system;
evaluating said second user's voice input with data values from said specified vocabulary according to an evaluation criterion; and 20
selecting a particular data value as an input into said computerised form field if said second user's voice input meets said criterion. 25

14. The method of claim 13 wherein said first and second users of said speech recognition system are different doctors and said computerised form fields are associated with a field within a computerised medical form. 30

15. A method as claimed in claim 1 comprising:

loading a second specified vocabulary from a second sub-database into computer storage, said second specified vocabulary associated with a second context used within said speech recognition system; 35
accepting said user's further voice input into said speech recognition system; 40
evaluating said user's voice input with data values from said specified vocabulary according to an evaluation criterion; and
selecting a particular data value as an input into said computerised form field if said user's voice input meets said evaluation criterion. 45

16. The method of claim 15 wherein said first context is a patient's age and said second context is a patient diagnosis of said patient. 50

17. A speech recognition system including a base vocabulary, the system comprising: 55

processing means (240,245) adapted to create a specified database containing text strings, wherein the text strings are provided from the

inputs of previous use of the system, and to define a sub-database within the specified database containing text strings associated with a context of input data;
a context identification module (240) adapted to identify the context of an input of data;
the processing means being further adapted to load (705) a specified vocabulary from the sub-database into computer storage, said specified vocabulary associated with a specific context;
to accepting (710) a user's voice input into said speech recognition system;
to evaluate (740) said user's voice input with data values from said specified vocabulary according to an evaluation criterion;
to select a particular data value as an input into a computerised form field if said evaluation criterion is met; and
if said user's voice input does not meet said evaluation criterion, to select a data value from said base vocabulary as an input into said computerised form field.

18. The speech recognition system of claim 17 wherein said context is a topical context.

19. The speech recognition system of claim 17 wherein said context is associated with a specific user of said speech recognition system.

20. The speech recognition system of claim 17 wherein said context is associated with said field.

21. A database for a speech recognition system, the database containing text strings, wherein the text strings are provided from the inputs of previous use of the system, the database comprising a plurality of sub-databases, each containing a respective vocabulary associated with a respective context of input data.

22. A computer program comprising instructions for controlling a processor of a speech recognition system when executed to carry out all of the steps of a method as claimed in any one of claims 1 to 16.

Patentansprüche

1. Verfahren zum Betrieb eines Spracherkennungssystems mit einem Basisvokabular, wobei in dem Verfahren eine spezifizierte Datenbank mit Textelementen erzeugt wird, wobei die Textelemente aus den Eingaben einer vorherigen Benutzung des Systems bereitgestellt werden, eine Sub-Datenbank innerhalb der spezifizierten Datenbank definiert wird, die Textelemente enthält, die

27

EP 1 599 867 B1

28

- einem Kontext von Eingabedaten zugeordnet sind, der Kontext von eingegebenen Daten identifiziert wird,
 ein spezifiziertes Vokabular aus der Sub-Datenbank in einen Computerspeicher geladen (705) wird, wobei das spezifizierte Vokabular einem spezifischen Kontext zugeordnet ist,
 eine Spracheingabe eines Benutzers in das Spracherkennungssystem aufgenommen (710) wird,
 die Spracheingabe des Benutzers mittels Datenwerten von dem spezifizierten Vokabular gemäß einem Bewertungskriterium bewertet (740) wird,
 ein bestimmter Datenwert als Eingabe in ein computerisiertes Formularfeld ausgewählt wird, falls das Bewertungskriterium erfüllt ist, und
 falls die Spracheingabe des Benutzers das Bewertungskriterium nicht erfüllt, ein Datenwert von dem Basisvokabular als eine Eingabe in das computerisierte Formularfeld ausgewählt wird.
2. Spracherkennungsverfahren nach Anspruch 1, wobei der Kontext dem Feld zugeordnet ist.
 3. Spracherkennungsverfahren nach Anspruch 1, wobei der Kontext einem thematischen Betreff zugeordnet ist.
 4. Spracherkennungsverfahren nach Anspruch 1, wobei der Kontext einem spezifischen Benutzer zugeordnet ist.
 5. Verfahren nach Anspruch 1, wobei die Spracheingabe des Benutzers ferner mittels Datenwerten aus dem Basisvokabular gemäß einem Basisbewertungskriterium bewertet wird, falls die Spracheingabe des Benutzers das Bewertungskriterium nicht erfüllt.
 6. Verfahren nach Anspruch 1, wobei das Bewertungskriterium in einer den Datenwerten zugeordneten Benutzungsgewichtung besteht.
 7. Verfahren nach Anspruch 1, wobei beim Bewertungsschritt ferner eine Übereinstimmungsheuristik gemäß einem bekannten Schwellenwert angewandt wird.
 8. Verfahren nach Anspruch 7, wobei beim Anwenden der Übereinstimmungsheuristik ferner die Spracheingabe des Benutzers mit einer Granzwahrscheinlichkeit verglichen wird, mit einem aus dem spezifizierten Vokabular abgeleiteten akustischen Modell übereinzustimmen.
 9. Verfahren nach Anspruch 1, wobei das erste spezifizierte Vokabular einem ersten computerisierten Formularfeld zugeordnet ist,
 ein zweites spezifiziertes Vokabular aus einer zweiten Sub-Datenbank in den Computerspeicher geladen wird, wobei das zweite spezifizierte Vokabular einem zweiten computerisierten Formularfeld zugeordnet ist,
 eine weitere Spracheingabe eines Benutzers in das Spracherkennungssystem aufgenommen wird,
 die Spracheingabe des Benutzers mittels Datenwerten von dem spezifizierten Vokabular gemäß einem Bewertungskriterium bewertet wird, und
 ein bestimmter Datenwert als Eingabe in ein zweites computerisiertes Formularfeld ausgewählt wird, falls die Spracheingabe des Benutzers das Bewertungskriterium erfüllt.
 10. Verfahren nach Anspruch 9, wobei die Bewertungskriterien für die Bewertungsschritte bezüglich des ersten und des zweiten spezifizierten Vokabulars gleich sind.
 11. Verfahren nach Anspruch 9, wobei die Bewertungskriterien für die Bewertungsschritte bezüglich des ersten und des zweiten spezifizierten Vokabulars verschieden voneinander sind.
 12. Verfahren nach Anspruch 9, wobei das erste und das zweite computerisierte Formularfeld verschiedenen Feldern eines computerisierten medizinischen Formulars zugeordnet sind.
 13. Verfahren nach Anspruch 1, wobei ein zweites spezifiziertes Vokabular aus einer zweiten Sub-Datenbank in einen Computerspeicher geladen wird, wobei das zweite spezifizierte Vokabular einem zweiten Benutzer des Spracherkennungssystems zugeordnet ist,
 eine Spracheingabe des zweiten Benutzers in das Spracherkennungssystem aufgenommen wird,
 die Spracheingabe des zweiten Benutzers mittels Datenwerten aus dem spezifizierten Vokabular gemäß einem Bewertungskriterium bewertet wird, und
 ein bestimmter Datenwert als eine Eingabe in das computerisierte Formularfeld ausgewählt wird, falls die Spracheingabe des zweiten Benutzers das Kriterium erfüllt.
 14. Verfahren nach Anspruch 13, wobei der erste und der zweite Benutzer des Spracherkennungssystems verschiedene Ärzte sind und die computerisierten Formularfelder einem Feld innerhalb eines computerisierten medizinischen Formulars zugeordnet sind.
 15. Verfahren nach Anspruch 1, wobei ein zweites spezifiziertes Vokabular aus einer zweiten Sub-Datenbank in den Computerspeicher geladen wird, wobei das zweite spezifizierte Vokabular einem zweiten innerhalb des Spracherkennungssystems benutzten Kontext zugeordnet ist,

29

EP 1 599 867 B1

30

- eine weitere Spracheingabe des Benutzers in das Spracherkennungssystem aufgenommen wird, die Spracheingabe des Benutzers mittels Datenwerten aus dem spezifizierten Vokabular gemäß einem Bewertungskriterium bewertet wird, und ein bestimmter Datenwert als eine Eingabe in das computerisierte Formularfeld ausgewählt wird, falls die Spracheingabe des Benutzers das Bewertungskriterium erfüllt.
16. Verfahren nach Anspruch 15, wobei der erste Kontext das Alter eines Patienten und der zweite Kontext eine Diagnose des Patienten ist.
17. Spracherkennungssystem mit einem Basisvokabular, mit einer Verarbeitungseinrichtung (240, 245), die dazu ausgelegt ist, eine spezifizierte Datenbank mit Textelementen, die aus den Eingaben einer vorherigen Benutzung des Systems bereitgestellt sind, zu erzeugen und eine Sub-Datenbank innerhalb der spezifizierten Datenbank zu definieren, die einem Kontext der Eingabedaten zugeordnete Textelemente enthält, einem Kontextidentifikationsmodul (240), das dazu ausgelegt ist, den Kontext von Eingabedaten zu identifizieren, wobei die Verarbeitungseinrichtung ferner dazu ausgelegt ist, ein spezifiziertes Vokabular aus der Sub-Datenbank in einen Computerspeicher zu laden (705), wobei das spezifizierte Vokabular einem spezifischen Kontext zugeordnet ist, eine Spracheingabe eines Benutzers in das Spracherkennungssystem aufzunehmen (710), die Spracheingabe des Benutzers mittels Datenwerten aus dem spezifizierten Vokabular gemäß einem Bewertungskriterium zu bewerten (740), einen bestimmten Datenwert als eine Eingabe in ein computerisiertes Formularfeld auszuwählen, falls das Bewertungskriterium erfüllt ist, und falls die Spracheingabe des Benutzers das Bewertungskriterium nicht erfüllt, einen Datenwert aus dem Basisvokabular als Eingabe in das computerisierte Formularfeld auszuwählen.
18. Spracherkennungssystem nach Anspruch 17, wobei der Kontext ein thematischer Kontext ist.
19. Spracherkennungssystem nach Anspruch 17, wobei der Kontext einem spezifischen Benutzer des Spracherkennungssystems zugeordnet ist.
20. Spracherkennungssystem nach Anspruch 17, wobei der Kontext dem Feld zugeordnet ist.
21. Datenbank für ein Spracherkennungssystem, wobei die Datenbank Textelemente enthält, die aus den

Eingaben einer vorherigen Benutzung des Systems bereitgestellt sind, wobei die Datenbank mehrere Sub-Datenbanken aufweist, die jeweils ein Vokabular enthalten, das einem entsprechenden Kontext von Eingabedaten zugeordnet ist.

22. Computerprogramm mit Befehlen zum Steuern eines Prozessors eines Spracherkennungssystems, um bei Ausführung des Programms alle Schritte eines Verfahrens gemäß einem der Ansprüche 1 bis 16 durchzuführen.

Revendications

1. Procédé de mise en oeuvre d'un système de reconnaissance vocale, ledit système de reconnaissance vocale comprenant un vocabulaire de base, le procédé comprenant :

la création d'une base de données spécifiée contenant des chaînes de texte, les chaînes de texte étant fournies à partir des entrées d'une utilisation précédente du système ;
la définition d'une sous-base de données dans la base de données spécifiée, contenant des chaînes de texte associées à un contexte de données d'entrée ;
l'identification du contexte d'une entrée de données ;
le chargement (705) d'un vocabulaire spécifié depuis la sous-base de données dans une mémoire d'ordinateur, ledit vocabulaire spécifié étant associé à un contexte spécifique ;
l'acceptation (710) d'une entrée vocale d'utilisateur dans ledit système de reconnaissance vocale ;
l'évaluation (740) de ladite entrée vocale d'utilisateur avec des valeurs de données provenant dudit vocabulaire spécifié conformément à un critère d'évaluation ;
la sélection d'une valeur de donnée particulière en tant qu'entrée dans un champ de forme informatisée si ledit critère d'évaluation est satisfait ; et
si ladite entrée vocale d'utilisateur ne satisfait pas auxdits critères d'évaluation, la sélection d'une valeur de donnée à partir dudit vocabulaire de base en tant qu'entrée dans ledit champ de forme informatisée.

2. Procédé de reconnaissance vocale selon la revendication 1, dans lequel ledit contexte est associé audit champ.
3. Procédé de reconnaissance vocale selon la revendication 1, dans lequel ledit contexte est associé à un sujet topique.

31

EP 1 599 867 B1

32

4. Procédé de reconnaissance vocale selon la revendication 1, dans lequel ledit contexte est associé à un utilisateur spécifique.
5. Procédé selon la revendication 1, comprenant en outre l'évaluation de ladite entrée vocale d'utilisateur avec des valeurs de données provenant dudit vocabulaire de base conformément à un critère d'évaluation de base si ladite entrée vocale d'utilisateur ne satisfait pas audit critère d'évaluation.
6. Procédé selon la revendication 1, dans lequel ledit critère d'évaluation est une pondération d'utilisation associée auxdites valeurs de données.
7. Procédé selon la revendication 1, dans lequel ladite étape d'évaluation comprend en outre l'étape d'application d'une heuristique d'adaptation par rapport à un seuil connu.
8. Procédé selon la revendication 7, dans lequel ladite étape d'application d'une heuristique d'adaptation comprend en outre une étape de comparaison de ladite entrée vocale d'utilisateur à une probabilité de seuil d'adaptation d'un modèle acoustique dérivé dudit vocabulaire spécifié.
9. Procédé selon la revendication 1, dans lequel ledit premier vocabulaire spécifié est associé à un premier champ de forme informatisée ; le chargement d'un second vocabulaire spécifié depuis une seconde sous-base de données dans une mémoire d'ordinateur, ledit second vocabulaire spécifié étant associé à un second champ de forme informatisée ; l'acceptation d'une autre entrée vocale d'utilisateur dans ledit système de reconnaissance vocale ; l'évaluation de ladite entrée vocale d'utilisateur avec des valeurs de données provenant dudit vocabulaire spécifié conformément à un critère d'évaluation ; et la sélection d'une valeur de donnée particulière en tant qu'entrée dans un second champ de forme informatisée si ladite entrée vocale d'utilisateur satisfait audit critère d'évaluation.
10. Procédé selon la revendication 9, dans lequel ledit critère d'évaluation pour lesdites étapes d'évaluation desdits premier et second vocabulaires spécifiés est le même.
11. Procédé selon la revendication 9, dans lequel lesdits critères d'évaluation pour lesdites étapes d'évaluation desdits premier et second vocabulaires spécifiés sont des critères différents.
12. Procédé selon la revendication 9, dans lequel lesdits premier et second champs de forme informatisée sont associés à différents champs d'une forme médicale informatisée.
13. Procédé selon la revendication 1, comprenant :
- le chargement d'un second vocabulaire spécifié depuis une seconde sous-base de données dans une mémoire d'ordinateur, ledit second vocabulaire spécifié étant associé à un second utilisateur dudit système de reconnaissance vocale ; l'acceptation d'une seconde entrée vocale d'utilisateur dans ledit système de reconnaissance vocale ; l'évaluation de ladite seconde entrée vocale d'utilisateur avec des valeurs de données provenant dudit vocabulaire spécifié conformément à un critère d'évaluation ; et la sélection d'une valeur de donnée particulière en tant qu'entrée dans ledit champ de forme informatisée si ladite seconde entrée vocale d'utilisateur satisfait audit critère.
14. Procédé selon la revendication 13, dans lequel lesdits premier et second utilisateurs dudit système de reconnaissance vocale sont des médecins différents et lesdits champs de forme informatisée sont associés à un champ dans une forme médicale informatisée.
15. Procédé selon la revendication 1, comprenant :
- le chargement d'un second vocabulaire spécifié depuis une seconde sous-base de données dans une mémoire d'ordinateur, ledit second vocabulaire spécifié étant associé à un second contexte utilisé dans ledit système de reconnaissance vocale ; l'acceptation de ladite autre entrée vocale d'utilisateur dans ledit système de reconnaissance vocale ; l'évaluation de ladite entrée vocale d'utilisateur avec des valeurs de données provenant dudit vocabulaire spécifié conformément à un critère d'évaluation ; et la sélection d'une valeur de donnée particulière en tant qu'entrée dans ledit champ de forme informatisée si ladite entrée vocale d'utilisateur satisfait audit critère d'évaluation.
16. Procédé selon la revendication 15, dans lequel ledit premier contexte est l'âge d'un patient et ledit second contexte est un diagnostic dudit patient.
17. Système de reconnaissance vocale comprenant un vocabulaire de base, le système comportant :
- un moyen de traitement (240, 245) conçu pour créer une base de données spécifiée contenant

33

EP 1 599 867 B1

34

des chaînes de texte, les chaînes de texte étant produites à partir des entrées d'une utilisation précédente du système, et pour définir une sous-base de données dans la base de données spécifiée contenant des chaînes de texte associées à un contexte de donnée d'entrée ;
 un module (240) d'identification de contexte conçu pour identifier le contexte d'une entrée de donnée ;
 le moyen de traitement étant en outre conçu pour charger (705) un vocabulaire spécifié depuis la sous-base de données dans une mémoire d'ordinateur, ledit vocabulaire spécifié étant associé à un contexte spécifique ;
 pour accepter (710) une entrée vocale d'utilisateur dans ledit système de reconnaissance vocale ;
 pour évaluer (740) ladite entrée vocale d'utilisateur avec des valeurs de données provenant dudit vocabulaire spécifié conformément à un critère d'évaluation ;
 pour sélectionner une valeur de donnée particulière en tant qu'entrée dans un champ de forme informatisée si ledit critère d'évaluation est satisfait ; et
 si ladite entrée vocale d'utilisateur ne satisfait pas audit critère d'évaluation, pour sélectionner une valeur de donnée provenant dudit vocabulaire de base en tant qu'entrée dans ledit champ de forme informatisée.

18. Système de reconnaissance vocale selon la revendication 17, dans lequel ledit contexte est un contexte topique.

19. Système de reconnaissance vocale selon la revendication 17, dans lequel ledit contexte est associé à un utilisateur spécifique dudit système de reconnaissance vocale.

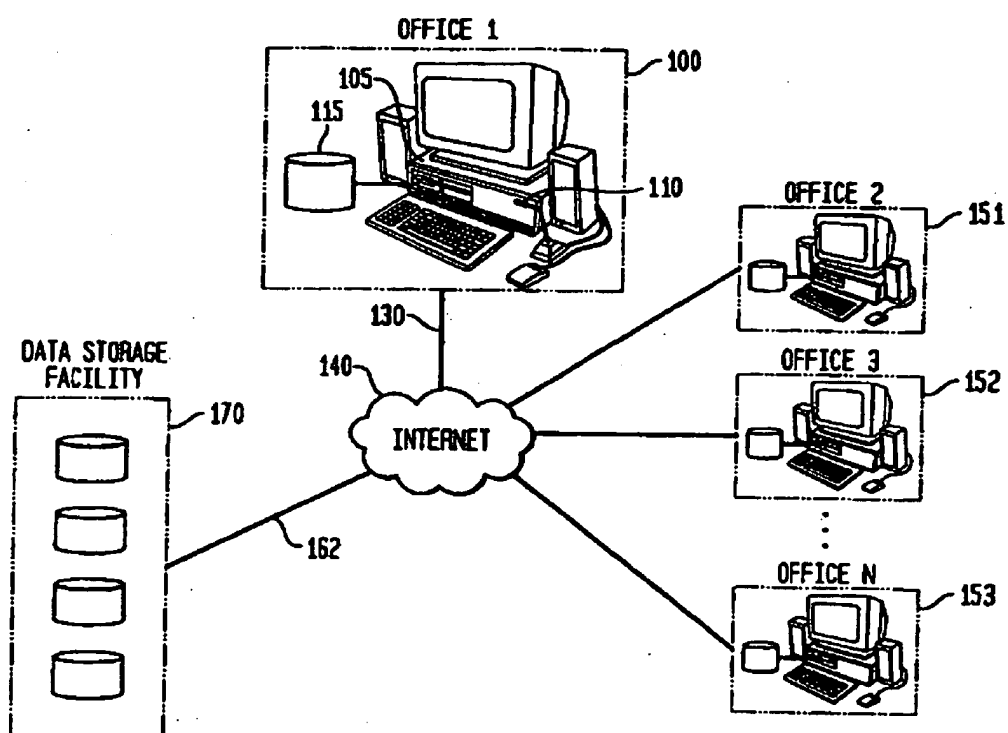
20. Système de reconnaissance vocale selon la revendication 17, dans lequel ledit contexte est associé audit champ.

21. Base de données pour un système de reconnaissance vocale, la base de données contenant des chaînes de texte, les chaînes de texte étant fournies depuis les entrées d'une utilisation précédente du système, la base de données comprenant de multiples sous-bases de données qui contiennent chacune un vocabulaire respectif associé à un contexte respectif de données d'entrée.

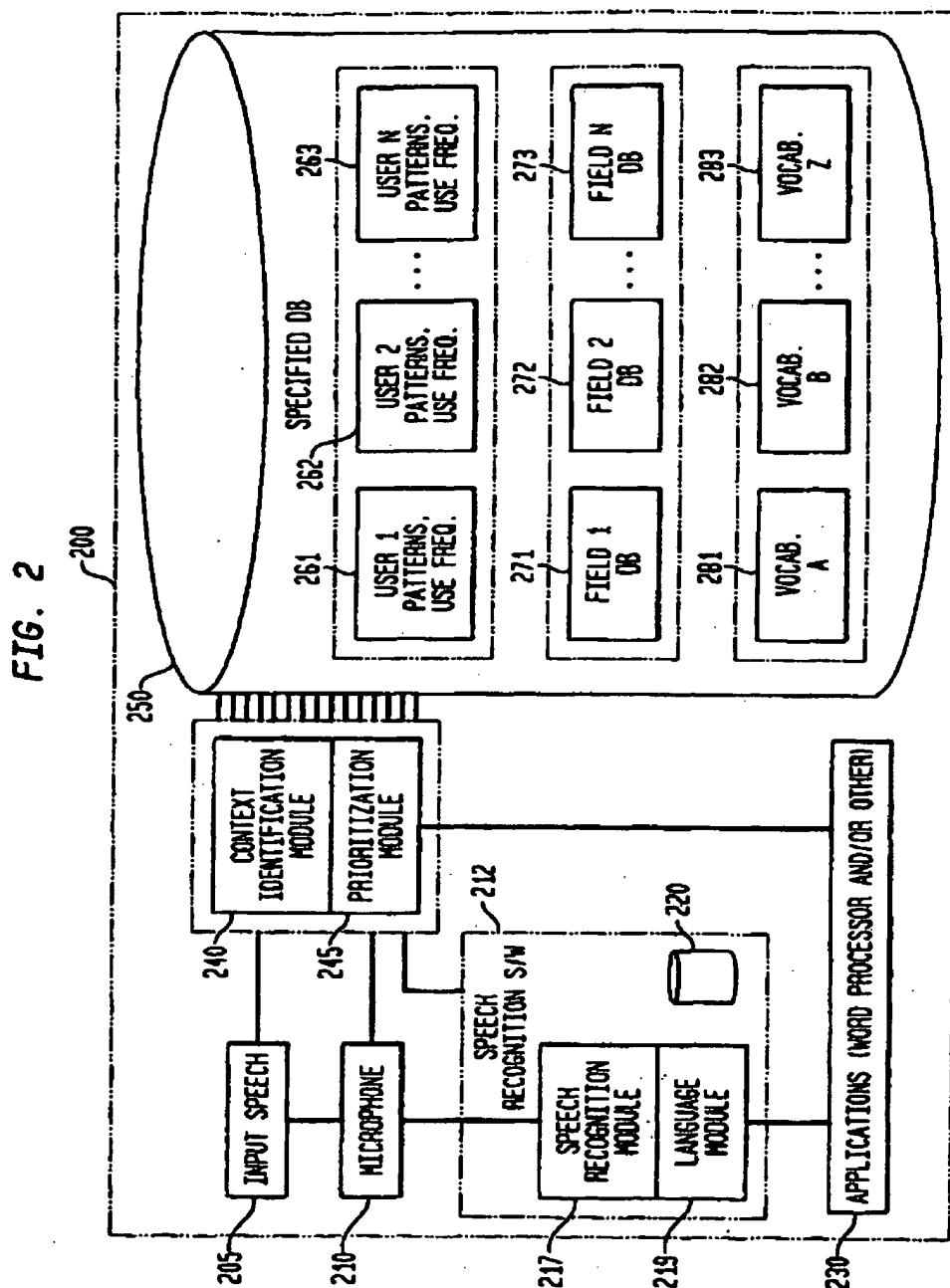
22. Programme d'ordinateur comportant des instructions pour commander un processeur d'un système de reconnaissance vocale lorsqu'il est mis en oeuvre pour exécuter toutes les étapes d'un procédé selon l'une quelconque des revendications 1 à 16.

EP 1 599 867 B1

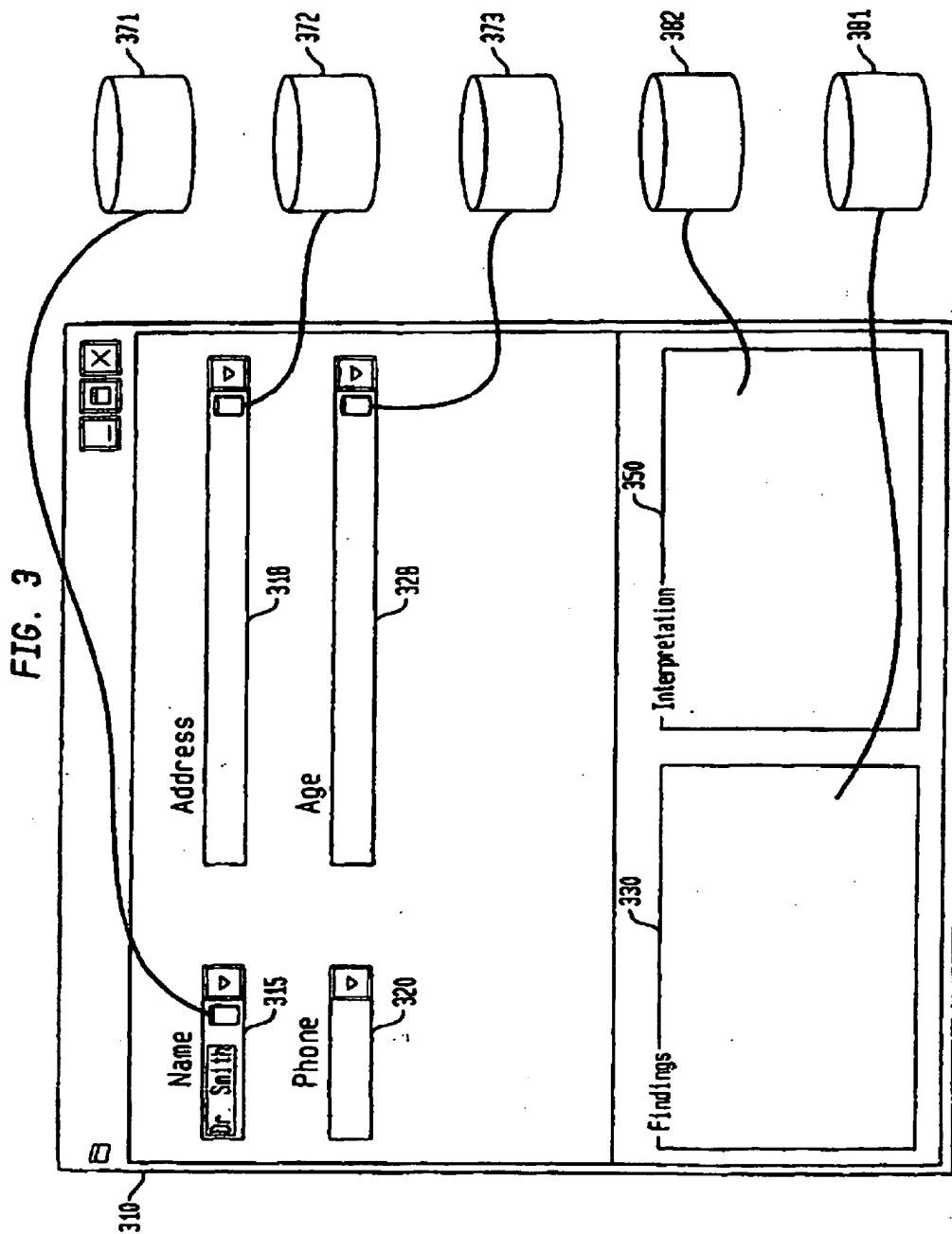
FIG. 1



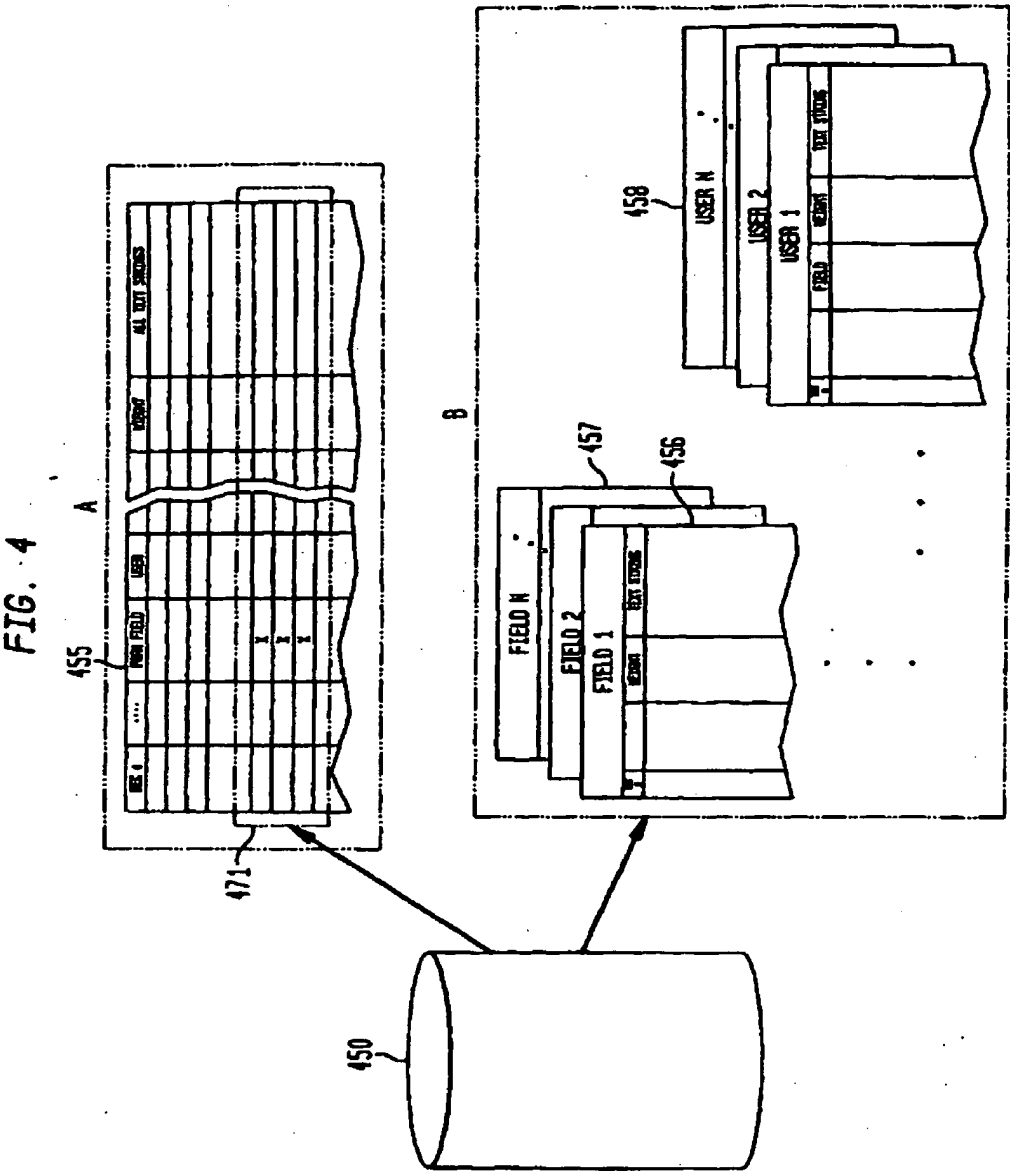
EP 1 599 867 B1



EP 1 599 867 B1



EP 1 599 867 B1



EP 1 599 867 B1

FIG. 5

554		552		571	
RECORD #		USERS	WEIGHTING INFO		STREET ADDRESS
5			7		10 WAYWARD LANE 510
7			10		15 BRANCH ST 511
65			2		555 5TH AVE 512

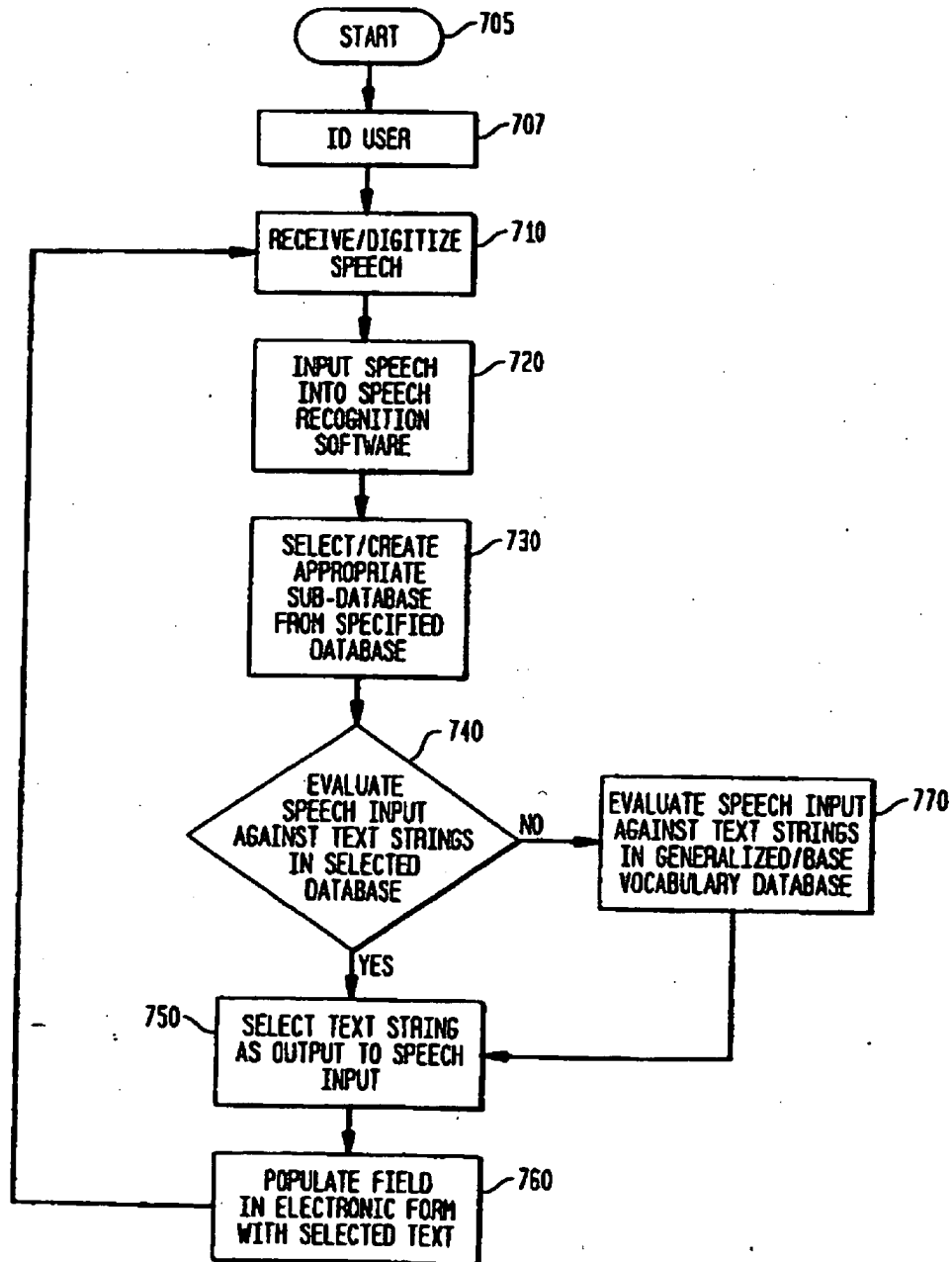
EP 1 599 867 B1

FIG. 6

RECORD #	CONTEXT (AGE)	CONTEXT (USER)	CONTEXT (FINDING)	WEIGHT (USE COUNT)	INTERPRETATIONS
3					PNEUMONIA
:					
:					
12	<3	DR. BROWN	SHALLOWING	5	DYSPHAGIA
13	>50	DR. SMITH	SPEECH	2	DYSPHAGIA

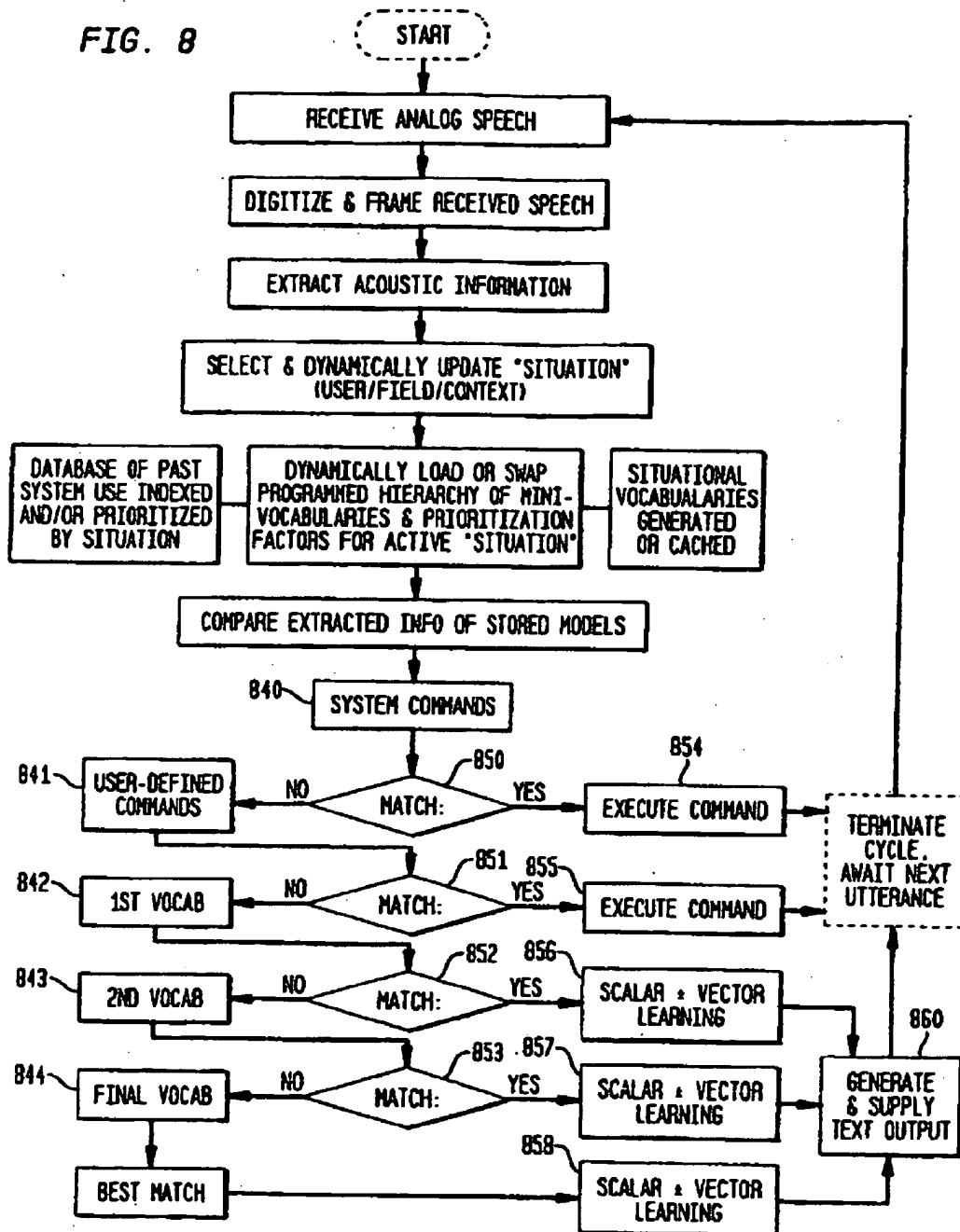
EP 1 599 867 B1

FIG. 7



EP 1 599 867 B1

FIG. 8



EP 1 599 867 B1**REFERENCES CITED IN THE DESCRIPTION**

This list of references cited by the applicant is for the reader's convenience only. It does not form part of the European patent document. Even though great care has been taken in compiling the references, errors or omissions cannot be excluded and the EPO disclaims all liability in this regard.

Patent documents cited in the description

- US 6631348 B, Wymore [0002]
- US 6662160 B, Chien [0003]
- US 6195837 B, Ballard [0004]
- US 6490557 B, Jeppesen [0005]
- US 6526380 B, Thelan [0006]
- WO 0126093 A [0008]
- WO 0189905 A [0009]
- US 6073097 A [0029]